

# A Constrained Conditional Likelihood Approach for Estimating the Means of Selected Populations

Claudio Fuentes<sup>1</sup> and Vik Gopal<sup>2</sup>

<sup>1</sup>*Department of Statistics, Oregon State University*

<sup>2</sup>*Department of Statistics and Applied Probability, National University of Singapore*

## Abstract

Given  $p$  independent normal populations, we consider the problem of estimating the mean of those populations, that based on the observed data, give the strongest signals. We explicitly condition on the ranking of the sample means, and consider a constrained conditional maximum likelihood (CCMLE) approach, avoiding the use of any priors and of any sparsity requirement between the population means. Our results show that if the observed means are too close together, we should in fact use the grand mean to estimate the mean of the population with the larger sample mean. If they are separated by more than a certain threshold, we should shrink the observed means towards each other. As intuition suggests, it is only if the observed means are far apart that we should conclude that the magnitude of separation and consequent ranking are not due to chance. Unlike other methods, our approach does not need to pre-specify the number of selected populations and the proposed CCMLE is able to perform simultaneous inference. Our method, which is conceptually straightforward, can be easily adapted to incorporate other selection criteria.

## 1 Introduction

Consider a scenario where  $p \geq 2$  independent normal populations are available, and from each one of them, we obtain a sample of size  $n$ . In this context, practitioners are sometimes interested in estimating the true means of those populations that yielded the  $k$  largest sample means in the experiment.

A naive solution to the problem is to estimate the selected population means with the corresponding sample means. Such an approach, however, is known to be problematic. [11] showed that the resulting estimator is biased for the case  $k = 1$ . This bias is particularly evident when all  $p$  populations are identically distributed. In terms of optimality, [16, 13] both showed that the estimator is minimax only when  $p = 2$ .

In order to improve on the naive estimator, several alternatives have been proposed in the literature, including [5], [2], and [3]. These papers propose estimators that perform better in terms of the Mean Squared Error (MSE). [18] considered a bias correction approach for the problem, obtaining estimators that perform well in terms of frequentist risk. Following up on his own idea, [19] introduced  $\omega$ -estimators which are essentially a weighted average of the order statistics. [4] considered a two-stage procedure, assuming we can obtain a second sample from the selected population, to produce unbiased estimators of the selected means. Despite these results, performance

theorems are scarce, with the exception of [1] and [9]. The latter proposes an empirical Bayes estimator and shows that it performs better in terms of the Bayes risk with respect to any normal prior. The former paper focuses on admissibility and obtains a generalized Bayes estimator under a harmonic prior.

More recently, [12] make an implicit assumption of sparsity, that many effect sizes (absolute population means)  $\theta_i = 0$ . By adapting the theory developed in [10], the authors in [12] perform post-selection inference with the Lasso. [15] approached the problem by estimating the first and second order bias of the naive estimator. Their results, which are similar in performance to the empirical Bayes approach in [6], extend to the non-Gaussian setting.

In this paper we propose a new estimator, that is based entirely on the likelihood function after incorporating the selection process. We motivate and define this new estimator in section 2, where we also provide a neat result for  $p = 2$  that yields some insight into how the estimator works. In section 3, we discuss the main computational hurdle in computing this estimator, and provide some direction on overcoming it. Following that section, we summarize the results of a simulation study that highlights the benefits and flaws of this new estimator. We conclude the paper with a brief discussion of our main results.

## 2 Conditional Likelihood Estimation

### 2.1 Defining the Conditional Likelihood

For simplicity, suppose that we obtain a single observation  $X_i \sim N(\mu_i, \sigma^2)$  from each population, and that the common variance  $\sigma^2 = 1$ . Then, for  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^T$ , the unconditional likelihood is given by  $L_0(\boldsymbol{\mu}) = \prod_{i=1}^p \phi(x_i - \mu_i)$ , where  $\phi$  denotes the density function of the standard normal distribution. However, once we observe  $X_1 = x_1, X_2 = x_2, \dots, X_p = x_p$ , we rank the observations in order to identify the populations corresponding to the largest  $x_i$ 's and estimate the respective means  $\mu_i$ . Hence, the conditional likelihood of interest is

$$L(\boldsymbol{\mu}) = \frac{\prod_{i=1}^p \phi(x_i - \mu_i)}{P_{\boldsymbol{\mu}}(X_1 > X_2 \dots > X_p)}, \quad (1)$$

where the notation  $P_{\boldsymbol{\mu}}(\cdot)$  explicitly states that the probability under consideration depends on  $\boldsymbol{\mu}$ . Note that, for equation (1), we do not have to worry about the labels attached to the groups. In other words, there is no loss of generality in assuming that the ordering of the means is  $X_1 > X_2 > \dots > X_p$ . It follows from equation (1) that the log of the conditional likelihood is

$$l(\boldsymbol{\mu}) = C - \frac{1}{2} \sum_{i=1}^p (x_i - \mu_i)^2 - \log P_{\boldsymbol{\mu}}(X_1 > X_2 \dots > X_p) \quad (2)$$

### 2.2 Constrained Conditional MLE

For the case  $p = 2$ , note that  $\boldsymbol{\mu} = (\mu_1, \mu_2)$  and

$$P_{\boldsymbol{\mu}}(X_1 > X_2) = P_{\boldsymbol{\mu}}(X_1 - X_2 > 0) = 1 - \Phi\left(\frac{\mu_2 - \mu_1}{\sqrt{2}}\right) \quad (3)$$

For a fixed  $\mu_2$ ,  $P_{\boldsymbol{\mu}}(X_1 > X_2)$  approaches 0 as  $\mu_1 \downarrow -\infty$ . This means that the  $-\log P_{\boldsymbol{\mu}}(X_1 > X_2)$  term in equation (2) increases to  $\infty$ , causing the conditional likelihood to be unbounded. In

other words, there is no global maximum for the expression in (2) for  $-\infty < \mu_1, \mu_2 < \infty$ . This phenomenon is of course, not specific to  $p = 2$ .

The occurrence of an unbounded likelihood is not without precedent in the statistics literature. Possibly the most widely studied models in which this occurs is a normal mixture model with unequal variances [7]. One of the solutions in that model was to constrain the parameter space to where there are local modes, and that is what we shall attempt to do here as well.

Since we expect that the populations with the largest sample means would be the ones with the largest population means, we now aim to maximise the conditional likelihood in equation (2) subject to  $\boldsymbol{\mu} \in \Theta$ , where  $\Theta = \{\boldsymbol{\mu} : \mu_1 \geq \mu_2 \geq \mu_3 \dots \geq \mu_p\}$ . We thus define the Constrained Conditional Maximum Likelihood Estimator (CCMLE) to be

$$\hat{\boldsymbol{\mu}} = \arg \max_{\boldsymbol{\mu} \in \Theta} l(\boldsymbol{\mu}) \quad (4)$$

This is clearly a case of constrained statistical inference. Much of the theory regarding this approach can be found in [14].

### 2.3 A Closed Form Result

**Proposition 1.** *Consider the CCMLE when  $p = 2$  and the variance is known. Thus, we have realisations  $x_1$  and  $x_2$  from  $N(\mu_i, \sigma^2)$  for  $i = 1, 2$ , and have observed that  $x_1 > x_2$ . If  $x_1 - x_2 > \frac{2\sigma}{\sqrt{\pi}}$ , then the unique CCMLE is an interior point of  $\Theta = \{(\mu_1, \mu_2) : \mu_1 \geq \mu_2\}$ . If  $x_1 - x_2 \leq \frac{2\sigma}{\sqrt{\pi}}$ , the CCMLE is  $\hat{\boldsymbol{\mu}} = (\bar{x}, \bar{x})$ , where  $\bar{x} = (x_1 + x_2)/2$ .*

*Proof.* Let us first define  $g(\cdot)$  to denote the inverse Mills ratio:

$$g(x) = \frac{\phi(x)}{1 - \Phi(x)} \quad (5)$$

We shall use certain properties of  $g$  in this proof – in particular, note that  $g$  is convex and monotone increasing [?, see]baricz2008mills.

Recall that the log-likelihood function is given by

$$l(\boldsymbol{\mu}) = -\log 2\pi - \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^2 (x_i - \mu_i)^2 - \log \left( 1 - \Phi \left( \frac{\mu_2 - \mu_1}{\sqrt{2\sigma^2}} \right) \right)$$

Setting the partial derivatives to zero, we find that, if a solution exists, it must satisfy

$$g \left( \sqrt{\frac{2}{\sigma^2}} (\bar{x} - \mu_1) \right) = \sqrt{\frac{2}{\sigma^2}} (x_1 - \mu_1) \quad (6)$$

Observe that equation (6) depends only on  $\mu_1$ . Let us define the function on the left to be  $h_1(\cdot)$  and the function on the right to be  $h_2(\cdot)$ . Thus, a stationary point exists if  $h_1$  and  $h_2$  intersect. Note that we are finding the root of a transcendental equation; no closed form solution exists. A plot of these two functions can be seen in Figure 1. Several quantities on the plot can be derived. For instance, from equation (5), we can compute that  $h_1(\bar{x}) = g(0) = \sqrt{\frac{2}{\pi}}$  and  $h_1'(\bar{x}) = -g'(0)\sqrt{2/\sigma^2} \approx -0.900/\sqrt{\sigma^2}$ . Moreover, the height  $H$  is  $(x_1 - x_2)/\sqrt{2\sigma^2}$ .

### Existence of Stationary Point

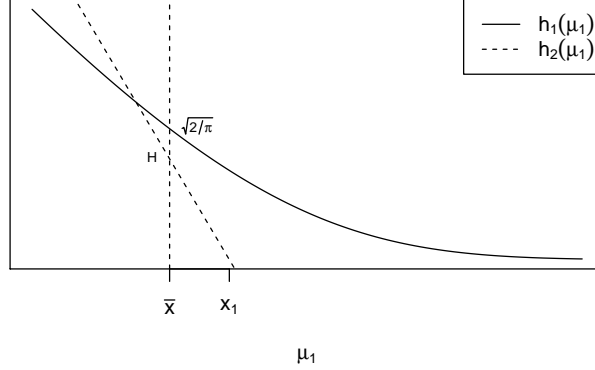


Figure 1: Conditions for existence/absence of a stationary point.

The figure shows that an intersection exists in  $\{\mu_1 : \mu_1 \geq \bar{x}\}$  if and only if  $H > \sqrt{2/\pi}$ . In other words, if and only if  $x_1 - x_2 > 2\sigma/\sqrt{\pi}$ .

Suppose that there is more than 1 root in  $\{\mu_1 : \mu_1 \geq \bar{x}\}$ . Then, at some point, it must be that  $h'_1(\mu_1) < -\sqrt{2/\sigma^2}$ . However  $h'_1(0) \approx -0.900/\sqrt{\sigma^2}$ , and  $h_1$  is convex. Hence  $h'_1(\mu_1) > -0.900/\sqrt{\sigma^2}$  for all  $\mu_1 \geq \bar{x}$ . It follows from this contradiction that the stationary point (if it exists) is unique.

Our final task is to show that, when  $x_1 - x_2 \leq 2\sigma/\sqrt{\pi}$ , the CCMLE is on the boundary of  $\Theta$  at  $(\bar{x}, \bar{x})$ . Within the constrained parameter space  $\Theta$ , we have  $\mu_2 \leq \mu_1$  and therefore

$$1 - \Phi\left(\frac{\mu_2 - \mu_1}{\sqrt{2\sigma^2}}\right) \geq 0.5$$

Hence, for all  $\boldsymbol{\mu} \in \Theta$ ,

$$l(\boldsymbol{\mu}) \leq C - \frac{1}{2\sigma^2} \sum_{i=1}^2 (x_i - \mu_i)^2 + \log 2$$

It is then easy to see that as  $\mu_1$  and/or  $\mu_2$  go to  $+\infty$  and/or  $-\infty$ , it holds that  $l(\boldsymbol{\mu}) \rightarrow -\infty$ .

When  $x_1 - x_2 < 2\sigma/\sqrt{\pi}$ , we already know that there are no stationary points in  $\Theta$ . This information, coupled with the understanding that  $l(\boldsymbol{\mu})$  decreases without bound as  $\boldsymbol{\mu}$  moves away from the boundary, leads to the conclusion that the CCMLE must be on the boundary  $\{\boldsymbol{\mu} : \mu_1 = \mu_2\}$ .

Now consider any point  $\boldsymbol{\mu}$  such that  $\mu_1 = \mu_2$ . If we can show that the directional derivative in the direction of  $(\bar{x}, \bar{x})$  is always positive, we are done, because it implies that we can always increase the log-likelihood by taking a suitable step in the direction of  $(\bar{x}, \bar{x})$ . To simplify notation, let

$$\gamma = \frac{1}{\sqrt{2\sigma^2}} g\left(\frac{\mu_2 - \mu_1}{\sqrt{2\sigma^2}}\right)$$

Thus we can write the gradient of  $l$  as  $\nabla l = \left(\frac{1}{\sigma^2}(x_1 - \mu_1) - \gamma, \frac{1}{\sigma^2}(x_2 - \mu_2) + \gamma\right)^T$ . The direction we need to consider is  $\mathbf{u} = K(\bar{x} - \mu_1, \bar{x} - \mu_2)^T$ , where  $K$  is a positive normalising constant, that makes  $|\mathbf{u}| = 1$ .

For any point in  $\Theta$ , the directional derivative is

$$\nabla l^T \mathbf{u} \propto ((1/\sigma^2)(x_1 - \mu_1) - \gamma)(\bar{x} - \mu_1) + ((1/\sigma^2)(x_2 - \mu_2) + \gamma)(\bar{x} - \mu_2) \quad (7)$$

Now for any point on the edge of  $\Theta$ , we let  $\mu = \mu_1 = \mu_2$  and simplify expression (7) to show that

$$\nabla l^T \mathbf{u} \propto \frac{2}{\sigma^2}(\bar{x} - \mu)^2 \geq 0$$

□

## 2.4 Sample Computations

It is not straightforward to generalize the proof in Proposition 1 to higher dimensions. However, the same intuitive results and numerical approaches apply.

In Table 1 below, we present some sample CCMLE computations for various combinations of observed values when  $p = 4$ . In performing the computations, we assume that  $\sigma^2$  is known, and is equal to 1. The purpose is to underline how, whenever a subset of observed sample means are close together, the CCMLE procedure will shrink them towards each other. This is apparent in row 1 of the table. Unlike the  $p = 2$  case, however, the critical distance at which they are collapsed onto one another is no longer  $2/\sqrt{\pi}$ . This latter phenomenon can be observed in row 2, where the separation is less than  $2/\sqrt{\pi}$ , but the estimated values are not exactly equal.

	Observed Values				Estimates			
	$x_1$	$x_2$	$x_3$	$x_4$	$\hat{\mu}_1$	$\hat{\mu}_2$	$\hat{\mu}_3$	$\hat{\mu}_4$
Config 1	10.0	9.5	9.0	0.0	9.50	9.50	9.50	0.00
Config 2	10.0	9.0	8.0	0.0	9.35	9.00	8.65	0.00
Config 3	10.0	9.0	1.0	0.0	9.50	9.50	0.50	0.50
Config 4	10.0	2.0	1.0	0.0	10.00	1.35	1.00	0.65

Table 1: Sample CCMLE computations for  $p = 4$

## 3 Computation of the CCMLE

### 3.1 Calculation of Probabilities

As part of numerically maximising the conditional likelihood in equation (2), we need to repeatedly compute  $P_{\boldsymbol{\mu}}(X_1 > X_2 \dots > X_p)$  for different  $\boldsymbol{\mu}$  vectors. In this section, we outline how it is possible to compute this integral for small values of  $p$  using a trick of conditioning, followed by a nested application of the `integrate` function in R.

For the case  $p = 2$ , it is clear from equation (3) that the probability of interest can be computed using the usual approximations to the standard normal distribution function.

For the case  $p = 3$ , we can condition on the random variable in the middle to yield a univariate integral.

$$\begin{aligned} P_{\boldsymbol{\mu}}(X_1 > X_2 > X_3) &= \int P_{\boldsymbol{\mu}}(X_1 > X_2 > X_3 | X_2 = x_2) f(x_2) dx_2 \\ &= \int P_{\boldsymbol{\mu}}(X_1 > x_2) P(x_2 > X_3) f(x_2) dx_2 \end{aligned}$$

The same approach enabled us to compute the probabilities up to  $p = 7$  without any further optimisation in R. For higher dimensions, we advocate computing the integral using a lower level language such as C, and switching to a sparse grid method [8] instead of persisting with cubature techniques.

### 3.2 Obtaining A Good Starting Point

We now focus on obtaining a good starting point for the numerical optimisation in the case that  $p > 2$ . A good starting point ensures that we can reduce the number of times that we evaluate the high-dimensional probability and its derivatives. Let us first denote  $f(\boldsymbol{\mu}) = \log P_{\boldsymbol{\mu}}(X_1 > X_2 \dots > X_p)$ . Taking a first order Taylor approximation of  $f$  about the observed  $\mathbf{x}$ , we can approximate the conditional likelihood in (2) with

$$l(\boldsymbol{\mu}) \approx C - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T(\mathbf{x} - \boldsymbol{\mu}) - f(\mathbf{x}) - (\boldsymbol{\mu} - \mathbf{x})^T \nabla f(\mathbf{x}) \quad (8)$$

Taking derivative with respect to  $\boldsymbol{\mu}$  and setting the above equation to 0, we can get an approximation to the CCMLE. It works out to be

$$\hat{\boldsymbol{\mu}}_0 \approx \mathbf{x} - \nabla f(\mathbf{x}) \quad (9)$$

The solution does not mean that a stationary point always exists - remember that we are solving an approximation to the conditional log-likelihood. For the same reason, it is also possible that the solution in equation (9) does not fall within the constrained space  $\Theta$ . In such cases, we shall use the orthogonal projection onto  $\Theta$  as the starting point. Note that the projection has to be performed numerically - fortunately however,  $\Theta$  is convex, and hence we can rely on prior methods from convex optimization theory in order to perform this step. The R package [17] provides a good implementation for solving this projection problem, using a quadratic programming technique.

## 4 A Comparison By Simulation

In this section, we shall assess the performance of the CCMLE via its Mean Squared Error, and the bootstrap confidence intervals constructed using it. We take the opportunity to highlight that the errors and intervals have to respect the selection procedure.

For instance, suppose that populations A, B and C have true means 3, 2 and 1 respectively. If the corresponding sample means are 2.1, 2.2 and 1.8, then the population selected as the maximum would be population B. The error in estimating the mean of the selected population using the sample mean would be

$$2 - 2.2 = -0.2$$

This is the methodology employed in Section 4.1. The error should not be computed as  $3 - 2.2 = +0.8$ , since the true mean of the selected population is in fact 2.

Similarly, when bootstrapping the strata in Section 4.2, the population selected to have the maximum will not always be population B. It could be population A or even population C depending on the bootstrap sample drawn. Thus it is not accurate to describe it as a confidence interval for population A (which has the largest mean). It is a confidence interval for the population *selected* to have the maximum mean.

#### 4.1 Mean Squared Error Comparison

In this subsection, we conduct a simulation study to understand the MSE of the CCMLE, as compared to the MSE of the ordinary MLE. We consider only the cases when  $p = 2$  and  $p = 3$  as they are sufficiently informative.

For the  $p = 2$  case, we fix  $\mu_2 = 0$ , and vary  $\mu_1$  from 0 to 5. For each configuration, we generate 1000 bivariate  $N(\boldsymbol{\mu}, \mathbf{I})$  random vectors, and then estimate the mean of the population with the larger sample mean.

The MSE estimate for each configuration has been plotted and smoothed in Figure 2. Notice that the CCMLE performs very well when  $\mu_1 - \mu_2 \leq 1.5$ . Beyond that, it performs approximately 10% worse than the unconditional MLE until the difference in means becomes quite large. From such a point onwards, it will be the case that the sample means will be far enough apart to warrant no shrinkage at all. It seems reasonable to guess that when populations are close together, the CCMLE will be very beneficial; however, when the population means are in fact far apart, it is not the right estimator to use. We shall witness this in the  $p = 3$  case as well.

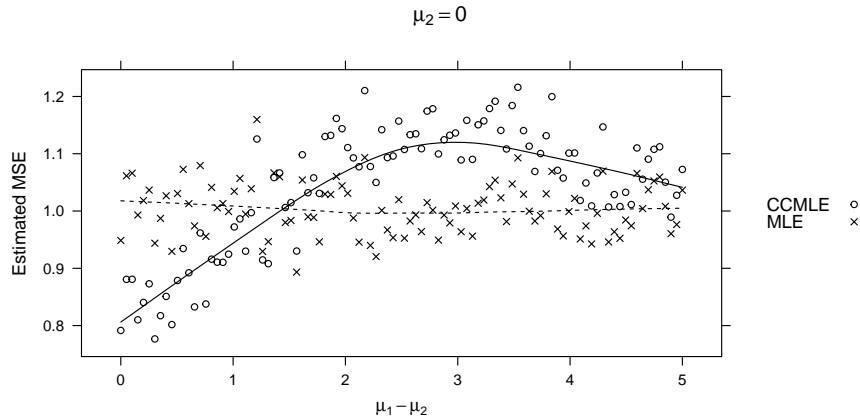


Figure 2: Estimated MSE for  $p = 2$  simulation experiment.

Now let us turn to the situation when  $p = 3$ . In our experiment, we considered 3 possible values for  $\mu_3$  : 0, 2 and 4. For each  $\mu_3$ , we varied  $\mu_2$  and  $\mu_1$  between  $\mu_3$  and 5. The output for these experiments is shown in Figure 3.

Consider the three separate 2-by-3 lattice plots in Figure 3. In each, the level plots in the top row correspond to the ordinary MLE, and the plots in the second row correspond to the CCMLE. The plots in the left-most column correspond to the errors when estimating the mean of the population with the *maximum* sample mean, and the plots in the right-most column correspond to the errors for the population with the *minimum* sample mean. Note that colors closer to black indicate a good estimator (low MSE) while colors closer to white indicate a poorer performance (high MSE).

Let us focus first on the case  $\mu_3 = 0$ . Observe that there is a diagonal dark line for the CCMLE in the left-most column. In the configurations on the diagonal, the means of population 1 and 2 are close to each other, and hence the CCMLE does well. In the off-diagonals, the populations are better separated and hence performance goes down. Overall though, it does appear that the MLE and CCMLE have a similar performance for this setting.

In the second configuration for  $p = 3$ , we had fixed  $\mu_3$  to be 2, while  $\mu_1$  and  $\mu_2$  varied from

2 to 5. The level plots for this configuration dark regions for the CCMLE than the MLE, when compared to the case when  $\mu_3 = 0$ . The reason is that the number of configurations where at least two of the populations are close together has increased, resulting in overall better performance from the CCMLE.

Now focus on the displayed plots for  $\mu_3 = 4$ , in Figure 3. The difference between the CCMLE and the MLE is more pronounced when we consider this final configuration, where  $\mu_1$  and  $\mu_2$  vary between 4 and 5. The regions in the level plot for the CCMLE are consistently darker than those for the MLE, due to the close proximity of sample means generated.

## 4.2 Confidence Intervals based on the CCMLE

In this subsection, we use the stratified bootstrap to assess the confidence intervals from the CCMLE procedure. Consider  $p = 3$ , and the true means to be  $\mu_1, \mu_2, \mu_3$ . We simulate a single sample of size 50 from each population. We draw from  $N(\mu_i, \sqrt{50})$  so that the sample mean still has variance 1. Then we draw 9999 bootstrap samples from within each sample and repeatedly compute the CCMLE. Each time, we are returned with an estimate of the mean of the populations selected to be the maximum, middle and minimum. The bias corrected intervals are then plotted for comparison with the traditional intervals that do not incorporate the selection process. This procedure is carried out for 4 different configurations of true means. The output can be seen in Figure 4.

The output is most informative for the configuration where  $\mu_1 = 10, \mu_2 = 9.5$  and  $\mu_3 = 9$ . In this case, the CCMLE point estimates are all equal, since there is insufficient power to discriminate between them. The CCMLE intervals are roughly the same, though it is worth pointing out the slight shrinkage, in opposite directions, of the intervals for the maximum and the minimum. In the situations where the mean of a group is far from the rest (for instance, in the bars for “Max” in the plots on the right, and the bars for “Min” in the bottom left), the CCMLE and the traditional approach provide comparable intervals. In cases where means are not distinguishable, (for instance, in the bars for “Mid” and “Min” in the top right), the CCMLE suggests it can provide shorter intervals.



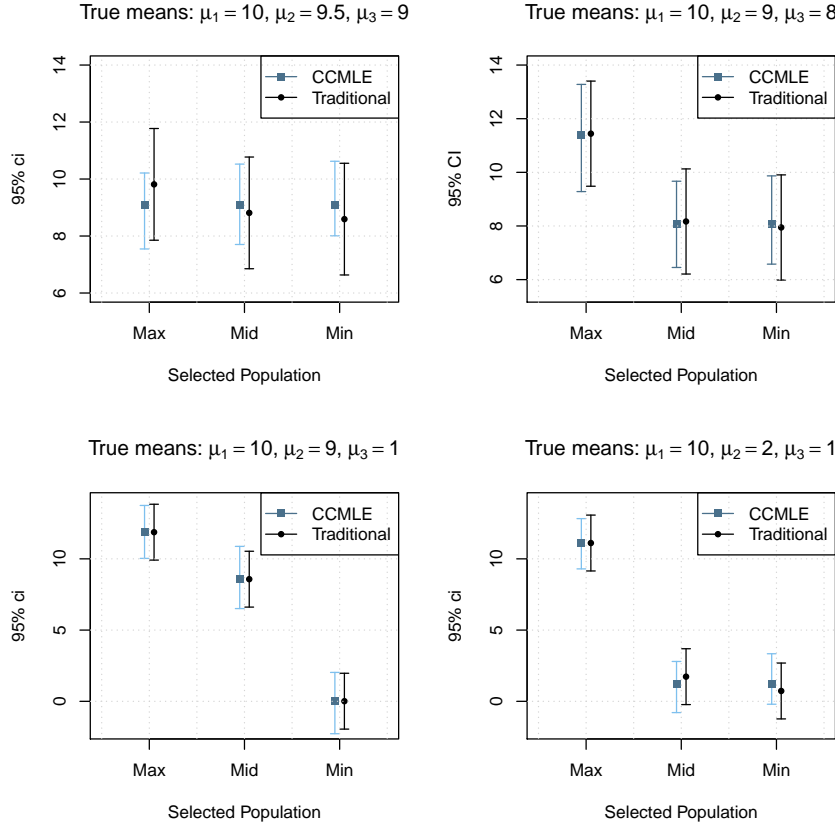


Figure 4: Bootstrap confidence intervals using CCMLE procedure. See Section 4.2 for a detailed explanation.

## 5 Discussion

In this paper, we introduced a new estimator for the means of selected populations and have started to unveil some of its properties. Although simulations experiments suggest the estimator is not admissible (see Figures 2 and 3), it performs well, particularly when the population means are close together. Furthermore, the proposed CCMLE provides simultaneous inference on selected populations, without a need to pre-specify the number of populations selected. Another advantage of the procedure is that because it is frequentist in paradigm, there is no need for prior specification.

Although the focus in this paper has been on selection via ranking of sample means, conceptually, this approach allows for any other selection criterion to be used. For instance, if populations were to be selected based on the absolute values of the sample means, the primary modification would be to the probability in the denominator of equation (1). Finally, as presented, it is straightforward to obtain bootstrap confidence intervals for this estimator, as demonstrated in Section 4.2.

Future work include the study of the asymptotic properties of this estimator, and computationally efficient methods to approximate the probability  $P_{\mu}(X_1 > \dots > X_p)$  in a general framework.

## References

- [1] Alexandra Bolotskikh. *Post-Selection Inference*. PhD thesis, Cornell University, 2015.
- [2] A. Cohen and H.B. Sackrowitz. Estimating the Mean of the Selected Population. In *Third Purdue Symposium on Statistical Decision Theory and Related Topics*. New York: Academic Press, 1982.
- [3] A. Cohen and H.B. Sackrowitz. A Decision Theoretic Formulation for Population Selection Followed by Estimating the Mean of the Selected Population. In *Fourth Purdue Symposium on Statistical Decision Theory and Related Topics*. New York: Academic Press, 1986.
- [4] Arthur Cohen and Harold B Sackrowitz. Two stage conditionally unbiased estimators of the selected mean. *Statistics & Probability Letters*, 8(3):273–278, 1989.
- [5] R.C. Dahiya. Estimation of the Mean of the Selected Population. *Journal of the American Statistical Association*, 69(345):226–230, 1974.
- [6] Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- [7] Richard J Hathaway. A constrained formulation of maximum-likelihood estimation for normal mixture distributions. *The Annals of Statistics*, pages 795–800, 1985.
- [8] Florian Heiss and Viktor Winschel. Likelihood approximation by numerical integration on sparse grids. *Journal of Econometrics*, 144(1):62–80, 2008.
- [9] J.T. Hwang. Empirical Bayes Estimation for the Means of the Selected Populations. *Sankhyā: The Indian Journal of Statistics, Series A*, 55(2):285–304, 1993.
- [10] Jason D Lee, Dennis L Sun, Yuekai Sun, and Jonathan E Taylor. Exact post-selection inference with applications to the lasso. *arXiv preprint arXiv:1311.6238*, 2014.
- [11] J. Putter and D. Rubinstein. On Estimating the Mean of a Selected Population. Technical Report 165, Department of Statistics, University of Wisconsin., 1968.
- [12] Stephen Reid, Jonathan Taylor, and Robert Tibshirani. Post-selection point and interval estimation of signal sizes in gaussian samples. *arXiv preprint arXiv:1405.3340*, 2014.
- [13] H. Sackrowitz and E. Samuel-Cahn. Evaluating the Chosen Population: A Bayes and Minimax Approach. *Lecture Notes-Monograph Series*, pages 386–399, 1986.
- [14] Mervyn J Silvapulle and Pranab Kumar Sen. *Constrained Statistical Inference: Order, Inequality, and Shape Constraints*, volume 912. John Wiley & Sons, 2011.
- [15] Noah Simon and Richard Simon. On estimating many means, selection bias, and the bootstrap. *arXiv preprint arXiv:1311.3709*, 2013.
- [16] C. Stein. Contribution to the Discussion of Bayesian and Non-Bayesian Decision Theory. *Handout from the Institute of Mathematical Statistics Meeting*, 1964.

- [17] Ravi Varadhan and Paul Gilbert. BB: An R package for solving a large system of nonlinear equations and for optimizing a high-dimensional nonlinear objective function. *Journal of Statistical Software*, 32(4):1–26, 2009.
- [18] JH Venter. Estimation of the Mean of the Selected Population. *Communications in Statistics-Theory and Methods*, 17(3):791–805, 1988.
- [19] JH Venter and SJ Steel. Estimation of the Mean of the Population Selected from k Populations. *Journal of Statistical Computation and Simulation*, 38(1-4):1–14, 1991.

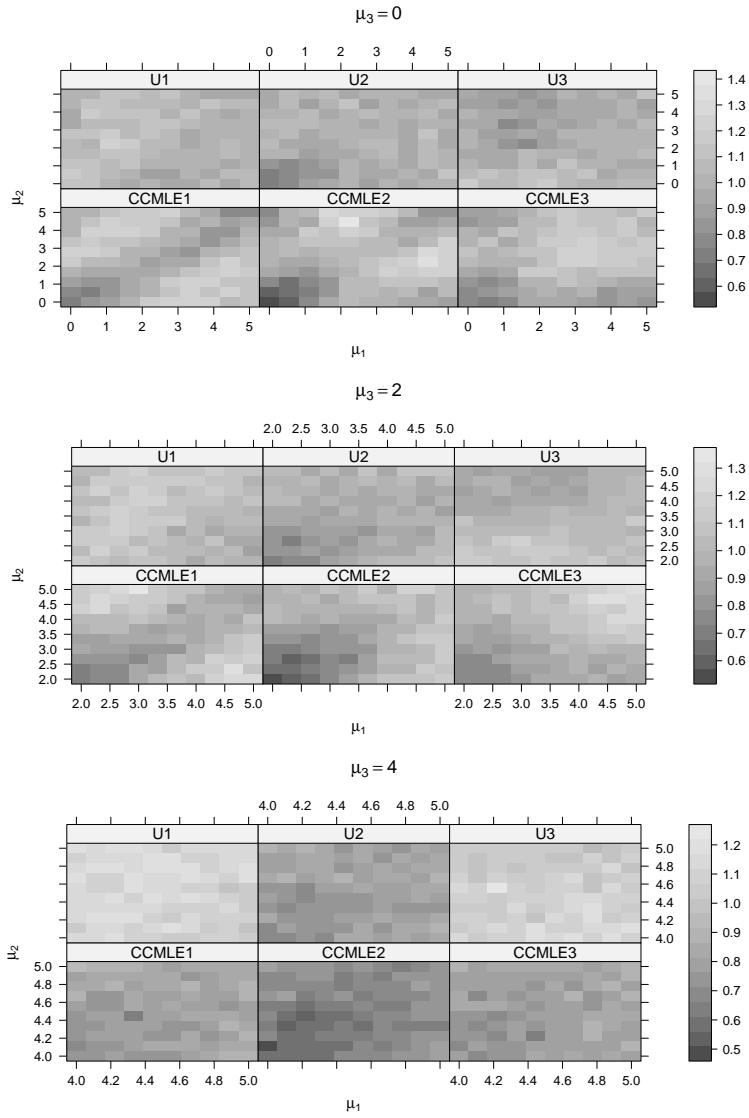


Figure 3: Estimated MSE for  $p = 3$  simulation experiment